

GHAJAR EXHIBIT 23

UNITED STATES DISTRICT COURT
NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO DIVISION

RICHARD KADREY, et al.,) Case No.
Individual and Representative) 3:23-cv-03417-VC
)
Plaintiffs,)
)
v.)
)
META PLATFORMS, INC., a)
Delaware corporation;)
)
Defendant.)
)

HIGHLY CONFIDENTIAL - ATTORNEYS' EYES ONLY
VIDEO-RECORDED DEPOSITION OF CRISTINA LOPES, PH.D.

Thursday, February 13, 2025
San Francisco, California

Stenographically Reported By:
Hanna Kim, CLR, CSR No. 13083
Job No. 7173634

1 THE COURT REPORTER: "Predict the next"?

2 THE WITNESS: -- the next token. So
3 it's -- it's con- -- the -- during training, it's
4 configured so that the inputs are the sequences of
5 text. The outputs are the same sequences of text, 10:03:44
6 but shifted by one, which gives the next over.

7 And so the -- the model goes through
8 several days, maybe weeks of doing that training
9 and -- until it converges to some acceptable
10 behavior. 10:04:05

11 BY MR. WEINSTEIN:

12 Q. At the end of a pretraining process, is
13 the result of that something called a base model?

14 A. That's how it's commonly called, yes.

15 Q. Okay. And a base model is a model that 10:04:15
16 has been pretrained but doesn't have any fine-tuning
17 applied to it; correct?

18 A. Correct.

19 Q. Okay. And in order to pretrain a large
20 language model, does that require a large amount of 10:04:31
21 data?

22 A. Yes.

23 Q. An enormous amount of data?

24 A. I mean, it depends on what you want to do
25 with the large language model. So there's -- 10:04:44

1 large -- you know, large language models can also be
2 for specific niche applications, like, for example,
3 code or software engineering tasks.

4 So it -- it depends if you -- you -- on
5 the other hand, if you want to train a large 10:04:55
6 language model to be multitasking,
7 multifunctionality, you know, do everything, then
8 you -- you -- you need a lot of data.

9 Q. All right.

10 Had you ever worked with -- yeah, 10:05:05
11 withdrawn.

12 Prior to your work on this case, had you
13 ever worked with any of the Meta large language
14 models?

15 A. No. 10:05:13

16 Q. Okay. But obviously since your work in
17 this case you've been working with some of the Llama
18 models; correct?

19 A. Yes.

20 Q. Okay. And you -- you have actually used 10:05:27
21 the Llama 3 7B and 70B models; correct?

22 A. Yes.

23 Q. Okay.

24 A. And 8B.

25 Q. 8B. 10:05:36

1 Q. Understood.

2 And that -- that's not just limited to
3 open-source data, that's any kind of text data --

4 A. Any --

5 Q. -- correct? 10:06:58

6 A. -- any kind of text data.

7 THE COURT REPORTER: We need to speak one
8 at a time, please.

9 THE WITNESS: I'm sorry.

10 MR. WEINSTEIN: I violated my own rule on 10:07:06
11 that.

12 BY MR. WEINSTEIN:

13 Q. Have you ever heard the term "overfitting"
14 as it refers to large language models?

15 A. Yes. 10:07:16

16 Q. What is "overfitting"?

17 A. It's when the network -- when the neural
18 network starts to learn too much about the data and
19 starts just basically memorizing the data and -- and
20 loses its ability to generalize from it. 10:07:30

21 Q. Is deduplication of training data an
22 example of a technique that helps prevent
23 overfitting?

24 A. Yes.

25 Q. Okay. What other techniques that -- 10:07:43

1 outputs and the obtained outputs. When that
2 happens, training can end, and the model is deemed
3 ready for use or 'inference.'" [As read]

4 That's correct; right?

5 A. Correct. 10:25:53

6 Q. And so just so we're clear, the term
7 "inference," that's a term that you use to describe
8 using the model after it's been at least in the
9 pretraining stage; correct?

10 A. Correct. 10:26:03

11 Q. Okay. And then Paragraph 40, you refer to
12 "generalization." And you say in the middle of the
13 Paragraph, "The value of deep learning models comes
14 from the value of mathematical regression itself:
15 the ability to produce reasonable outputs even when 10:26:27
16 the inputs were not in the training data. This is
17 called 'generalization.'"

18 Do you see that?

19 A. Yes.

20 Q. Is generalization a desired characteristic 10:26:33
21 of a large language model?

22 A. Yes.

23 Q. And people who build large language
24 models, they strive to make their models generalized
25 as much as possible; correct? 10:26:55

1 MR. YOUNG: Objection to form.

2 THE WITNESS: I believe so.

3 BY MR. WEINSTEIN:

4 Q. Okay. And by generalization, we're
5 referring to the ability of the model to provide 10:27:12
6 output that is different from the data on which it
7 was trained; correct?

8 A. It's -- it refers to the ability of the
9 model to produce reasonable outputs on data on which
10 it was not trained. 10:27:31

11 Q. Right.
12 Obviously, it's -- it's not any random
13 data, you want reasonable outputs; right?

14 A. Yes.

15 Q. The -- the goal is to provide an output 10:27:39
16 that was not in the original training data; correct?

17 A. Correct.

18 Q. Right.
19 You also say that, in Paragraph 40,
20 "Generalization stands on the quality of the 10:27:51
21 training data."

22 Do you see that?

23 A. Yes.

24 Q. When you talk about training data quality,
25 that's not necessarily the same thing as the quality 10:28:00

1 BY MR. WEINSTEIN:

2 Q. So by "high quality," is this referring to
3 the property of the data, that it's expected to
4 cause an improvement in the model's performance?

5 A. Yes. 10:30:43

6 Q. Okay. And then you mentioned Wikipedia.
7 And I realize you may not have researched
8 this, but what is the reason that Wikipedia is
9 regarded as high-quality training data, assuming it
10 is? 10:31:01

11 MR. YOUNG: Objection. Form.

12 Go ahead.

13 THE WITNESS: I can -- I -- I -- I don't
14 know. I can give an educated guess. It's because
15 it's highly curated. It has a lot of eyes on it. 10:31:10
16 The text tends to be very good. The explanations
17 tend to be very factual. There's -- references back
18 up the statements in most of the articles.

19 So I would imagine that that's the reason
20 why the LLM developers like Wikipedia so much. 10:31:30

21 BY MR. WEINSTEIN:

22 Q. And as -- as far as books you mentioned,
23 why, in your opinion, do books provide a high
24 quality of training data?

25 A. Again, I am trying -- I -- I -- I have -- 10:31:42

1 Common Crawl has 61 trillion tokens.

2 Do you see that?

3 A. Yes.

4 Q. You'd agree that not all the tokens in

5 Common Crawl are human readable text; correct? 10:45:36

6 A. I don't know.

7 Q. Well, for example, there may be HTML

8 elements in the web pages that are crawled

9 [verbatim]; correct?

10 A. I would imagine, yes. 10:45:50

11 Q. There may be nontext media elements;

12 correct?

13 A. Probably, yes.

14 Q. URLs?

15 A. Yes. 10:45:58

16 Q. Metadata?

17 A. Yes.

18 Q. Okay. You reference quality as being an
19 important consideration in training data.

20 Would you agree that's not the only 10:46:18
21 important consideration in terms of training data
22 used for an LLM?

23 A. Yes. So -- yeah.

24 Q. You also need a lot of training data;

25 correct? 10:46:33

1 A. That's correct.

2 Q. And you need a wide diversity of training
3 data; correct?

4 A. So, as I said, it depends on how you want
5 to use your large language model for; right? 10:46:44

6 So if you're talking about a -- a large
7 language model that wants to do everything, then you
8 want a lot of diversity. You want to be able to
9 capture everything.

10 Q. And Llama would be an example of a large 10:46:54
11 language model that would want to do everything?

12 MR. YOUNG: Objection to form.

13 THE WITNESS: I believe so.

14 BY MR. WEINSTEIN:

15 Q. The same would probably apply to ChatGPT; 10:47:06
16 correct?

17 A. Yeah.

18 MR. YOUNG: Object to form.

19 BY MR. WEINSTEIN:

20 Q. If we could turn to Paragraph 87, there's 10:47:15
21 a reference to "Meta downloaded a dataset known as
22 'The Pile.'" [As read]

23 Do you see that?

24 A. Mm-hmm. Yes.

25 Q. Now, I'll represent to you there's another 10:47:36

1 Q. Possibly thousands of times a day; is that
2 correct?

3 MR. YOUNG: Objection to form.

4 THE WITNESS: I don't know the number.

5 BY MR. WEINSTEIN: 02:17:30

6 Q. Okay. And so, would you agree with him
7 that the likely reason that his tests showed a
8 higher degree of memorization for Embraced was
9 because of the presence of Bible passages that were
10 otherwise available on the internet? 02:17:47

11 A. It's a -- it's a possible explanation,
12 yes.

13 Q. Okay. Why don't we turn to Paragraph 91.
14 And this is a looking at some tests that you ran
15 using Llama 3 70B; is that right? 02:18:37

16 A. Yes.

17 Q. And these were tests that were run on a
18 server at UCI?

19 A. Yes.

20 Q. Okay. Using a base model; correct? 02:18:48

21 A. Correct.

22 Q. Okay. And in this experiment, you said
23 that "Counsel provided me with 73 passage drawn from
24 a variety of books, including Asserted Works, to be
25 used as prompts for a completion using the Llama 8B 02:19:06

1 and the Llama 3 70B models." [As read]

2 Do you see that?

3 A. Yes.

4 Q. So your tests were using the same

5 technique of providing text from a book and seeing 02:19:16

6 if the model could complete it; correct?

7 A. Correct.

8 Q. Okay. So the 73 passages that you were

9 provided, they were selected by counsel for you?

10 A. Yes. 02:19:34

11 Q. Okay. And in a footnote, you're citing to

12 a -- a couple of LinkedIn posts, it looked like

13 from -- from the initials, Louis W. Hunt?

14 A. Yes.

15 Q. Have you ever spoken to Louis -- 02:19:54

16 A. No.

17 Q. -- Hunt?

18 A. No.

19 Q. Do you know who he is?

20 A. No. 02:20:00

21 Q. Okay. The samples that you were

22 provided -- withdrawn.

23 The passages that you were provided, were

24 they passages that were from Mr. Hunt's work?

25 A. Yes, I believe so. 02:20:18

1 Q. Okay. So Mr. Hunt published a series of
2 examples of passages and showed a certain degree of
3 completion for some of the passages using Llama; is
4 that right?

5 A. That's right. 02:20:40

6 Q. Okay. And the footnotes show in one case,
7 he had 400 pages of algorithmically generated
8 activity. The other 131, did you look at any of the
9 other ones that he did other than the 73 that
10 counsel provided to you? 02:21:01

11 A. No.

12 Q. Just the 73, okay.

13 A. Yes.

14 Q. What was the -- see if there's a way I can
15 ask this without getting into Rule 26 territory. 02:21:22

16 What was the basis for the selection of
17 the 73 passages to which you were provided?

18 MR. YOUNG: I'm going to caution the
19 witness to not reveal any discussions she may have
20 had with counsel or any -- reveal any mental 02:21:35
21 impressions she -- of counsel.

22 But to the extent you are able to answer,
23 please do.

24 THE WITNESS: Okay. Can you repeat the
25 question. 02:21:48

1 a large language model to be able to repeat part of
2 input text is a function of how many times the input
3 text was present in the dataset; correct?

4 A. Yes, I understand that.

5 Q. So would it be your expectation that when 03:32:44
6 Llama is able to complete certain passages, that
7 those passages likely appear multiple times in the
8 dataset?

9 A. Yes. That's why de-duplication is so
10 important. 03:33:02

11 Q. Now, all these tests that were done to
12 confirm the ability to complete passages, they all
13 involved inputting text from the original work as
14 the prompt; correct?

15 A. Yes. 03:33:19

16 Q. So to -- in -- in order to run these
17 tests, somebody who created these prompts would have
18 had to have access to the original text; correct?

19 A. Yes.

20 MR. YOUNG: Objection to form. 03:33:31

21 BY MR. WEINSTEIN:

22 Q. Have you heard the term "adversarial
23 prompt" before?

24 A. Yes.

25 Q. What's an adversarial prompt? 03:33:40

1 A. These are prompts that people use to try
2 to coerce these large language models to misbehave
3 in some way.

4 Q. Understood.

5 The prompts that were identified for 03:33:53
6 purposes of these memorization studies in Appendix C
7 of your report, are those examples of adversarial
8 prompts?

9 A. I don't think so. I think these are just
10 typical examples of, you know, trying to complete 03:34:07
11 things. Now, the -- the collection -- if you ask me
12 about the collection, like who collected these
13 prompts, it was not me. These prompts were given to
14 me. I -- I don't -- you know, I don't know why,
15 when, what went into the selection of these prompts, 03:34:29
16 as I said before.

17 Q. Right. But the per- -- person who put --
18 put together the prompts was trying to get the model
19 to output text from these books; correct?

20 MR. YOUNG: Objection to form and scope. 03:34:41

21 THE WITNESS: I -- I think that's what he
22 was trying to do. I mean, but that's also what all
23 of these memorization experiments are trying to do;
24 right?

25 BY MR. WEINSTEIN: 03:34:52

1 Q. Right. But this type of use of a large
2 language model to complete passages with original
3 text, is that a use you've done in -- is that part
4 of ordinary use of a large language model?

5 A. No. This is part of checking the -- you 03:35:09
6 know, the soundness and the -- the behavior of large
7 language models is part of something that -- you
8 know, like red teaming does, for example. So
9 this --

10 THE COURT REPORTER: Was the what? 03:35:19

11 THE WITNESS: Red teaming. So groups that
12 are, even inside the same development, the same
13 company or maybe in other companies that try to
14 really poke at the system and make it mis- --
15 misbehave so that, you know, measures and 03:35:38
16 mitigations can be taken.

17 BY MR. WEINSTEIN:

18 Q. And -- and all these tests were run on the
19 base model; correct?

20 A. Yes. 03:35:47

21 Q. So if there were mitigations in place to
22 try to prevent output that repeats training data as
23 part of post training, those wouldn't have been
24 effective in these experiments; correct?

25 MR. YOUNG: Form and scope. 03:36:02

CERTIFICATE OF REPORTER

I, Hanna Kim, a Certified Shorthand Reporter, do hereby certify:

That prior to being examined, the witness in the foregoing proceedings was by me duly sworn to testify to the truth, the whole truth, and nothing but the truth;

That said proceedings were taken before me at the time and place therein set forth and were taken down by me in shorthand and thereafter transcribed into typewriting under my direction and supervision;

I further certify that I am neither counsel for, nor related to, any party to said proceedings, not in anywise interested in the outcome thereof.

Further, that if the foregoing pertains to the original transcript of a deposition in a federal case, before completion of the proceedings, review of the transcript [X] was [] was not requested.

In witness whereof, I have hereunto subscribed my name.

Dated: February 28, 2025.



Hanna Kim, CLR, CSR No. 13083